

Approximate Maximum Likelihood Estimation for Linear Regression with Operator Uncertainty

Richard Clancy
(joint with Stephen Becker)

University of Colorado at Boulder
Department of Applied Mathematics

March 13th
FRAMSC 2021

Problem Formulation

Maximum Likelihood Estimation and Its Limitations

Moment Generating Functions and Saddle Point Approximation

Approximate Likelihood Function and Optimization Problem

Algorithm and Numerical Experiments

Table of Contents

Problem Formulation

Maximum Likelihood Estimation and Its Limitations

Moment Generating Functions and Saddle Point Approximation

Approximate Likelihood Function and Optimization Problem

Algorithm and Numerical Experiments

Problem Setup

We consider the generative model $\mathbf{y} = \mathbf{G}\mathbf{x} + \boldsymbol{\eta}$

- ▶ $\mathbf{G} \in \mathbb{R}^{m \times n}$ is a random matrix
- ▶ $\boldsymbol{\eta} \in \mathbb{R}^m$ is a random vector
- ▶ $\mathbf{x} \in \mathbb{R}^n$ is vector of model parameters
- ▶ $\mathbf{y} \in \mathbb{R}^m$ is a vector of measurements

Problem Setup

We consider the generative model $\mathbf{y} = \mathbf{G}\mathbf{x} + \boldsymbol{\eta}$

- ▶ $\mathbf{G} \in \mathbb{R}^{m \times n}$ is a random matrix
- ▶ $\boldsymbol{\eta} \in \mathbb{R}^m$ is a random vector
- ▶ $\mathbf{x} \in \mathbb{R}^n$ is vector of model parameters
- ▶ $\mathbf{y} \in \mathbb{R}^m$ is a vector of measurements

Given measurement vector \mathbf{y} and distributional knowledge of \mathbf{G} and $\boldsymbol{\eta}$, estimate \mathbf{x} .

Why should we care?

In practice, it is uncommon to know \mathbf{G} precisely. Some causes are

- ▶ precision limits in measurement
- ▶ truncation error for memory savings
- ▶ sampling error
- ▶ human error
- ▶ modeling error

Estimating home prices

Generative model $\mathbf{y} = \mathbf{G}\mathbf{x} + \boldsymbol{\eta}$

Estimating home prices

Generative model $\mathbf{y} = \mathbf{G}\mathbf{x} + \boldsymbol{\eta}$

Example

- ▶ $\mathbf{G} \in \mathbb{R}^{m \times 2}$, first column is square footage of home, second is square footage of lot. We observe $\mathbf{H} = \text{round}(\mathbf{G})$ to nearest hundred foot (e.g. 1366 \rightarrow 1400 sqft house)

Estimating home prices

Generative model $\mathbf{y} = \mathbf{G}\mathbf{x} + \boldsymbol{\eta}$

Example

- ▶ $\mathbf{G} \in \mathbb{R}^{m \times 2}$, first column is square footage of home, second is square footage of lot. We observe $\mathbf{H} = \text{round}(\mathbf{G})$ to nearest hundred foot (e.g. 1366 \rightarrow 1400 sqft house)
- ▶ G_{ij} can be modeled as a Uniform($H_{ij} - \delta, H_{ij} + \delta$) for $\delta = 50$ and $H_{ij} = 1400$

Estimating home prices

Generative model $\mathbf{y} = \mathbf{G}\mathbf{x} + \boldsymbol{\eta}$

Example

- ▶ $\mathbf{G} \in \mathbb{R}^{m \times 2}$, first column is square footage of home, second is square footage of lot. We observe $\mathbf{H} = \text{round}(\mathbf{G})$ to nearest hundred foot (e.g. 1366 \rightarrow 1400 sqft house)
- ▶ G_{ij} can be modeled as a Uniform($H_{ij} - \delta, H_{ij} + \delta$) for $\delta = 50$ and $H_{ij} = 1400$
- ▶ $\mathbf{y} \in \mathbb{R}^m$ is selling price for corresponding home (e.g. \$207k)

Estimating home prices

Generative model $\mathbf{y} = \mathbf{G}\mathbf{x} + \boldsymbol{\eta}$

Example

- ▶ $\mathbf{G} \in \mathbb{R}^{m \times 2}$, first column is square footage of home, second is square footage of lot. We observe $\mathbf{H} = \text{round}(\mathbf{G})$ to nearest hundred foot (e.g. 1366 \rightarrow 1400 sqft house)
- ▶ G_{ij} can be modeled as a Uniform($H_{ij} - \delta, H_{ij} + \delta$) for $\delta = 50$ and $H_{ij} = 1400$
- ▶ $\mathbf{y} \in \mathbb{R}^m$ is selling price for corresponding home (e.g. \$207k)
- ▶ GOAL: Estimate parameters \mathbf{x} so we can model price based on home and lot size accounting for uncertainty in \mathbf{G} and $\boldsymbol{\eta}$

Table of Contents

Problem Formulation

Maximum Likelihood Estimation and Its Limitations

Moment Generating Functions and Saddle Point Approximation

Approximate Likelihood Function and Optimization Problem

Algorithm and Numerical Experiments

Maximum likelihood estimation

- ▶ Given observed data \mathbf{y} and a likelihood function $L(\mathbf{x}) = f_{\mathbf{Y}}(\mathbf{y}; \mathbf{x})$, where $f_{\mathbf{Y}}(\mathbf{y}; \mathbf{x})$ is the PDF of \mathbf{y} , find parameters \mathbf{x} that maximize the likelihood function,

$$\operatorname{argmax}_{\mathbf{x}} L(\mathbf{x})$$

Maximum likelihood estimation

- ▶ Given observed data \mathbf{y} and a likelihood function $L(\mathbf{x}) = f_{\mathbf{Y}}(\mathbf{y}; \mathbf{x})$, where $f_{\mathbf{Y}}(\mathbf{y}; \mathbf{x})$ is the PDF of \mathbf{y} , find parameters \mathbf{x} that maximize the likelihood function,

$$\operatorname{argmax}_{\mathbf{x}} L(\mathbf{x})$$

- ▶ When components of \mathbf{y} are independent, we can split PDF such that

$$f_{\mathbf{Y}}(\mathbf{y}; \mathbf{x}) = \prod_{i=1}^m f_{Y_i}(y_i; \mathbf{x})$$

Maximum likelihood estimation

- ▶ Given observed data \mathbf{y} and a likelihood function $L(\mathbf{x}) = f_{\mathbf{Y}}(\mathbf{y}; \mathbf{x})$, where $f_{\mathbf{Y}}(\mathbf{y}; \mathbf{x})$ is the PDF of \mathbf{y} , find parameters \mathbf{x} that maximize the likelihood function,

$$\operatorname{argmax}_{\mathbf{x}} L(\mathbf{x})$$

- ▶ When components of \mathbf{y} are independent, we can split PDF such that

$$f_{\mathbf{Y}}(\mathbf{y}; \mathbf{x}) = \prod_{i=1}^m f_{Y_i}(y_i; \mathbf{x})$$

- ▶ We focus on maximizing the log-likelihood function

$$\hat{\mathbf{x}}_{MLE} = \operatorname{argmax}_{\mathbf{x}} \left\{ \sum_{i=1}^m \ln [f_{Y_i}(y_i; \mathbf{x})] \right\}$$

Justification of MLE for regression problem

For additive noise models $\mathbf{y} = \mathbf{H}\mathbf{x} + \boldsymbol{\eta}$ where \mathbf{H} is known precisely (all uncertainty is in $\boldsymbol{\eta}$), we have

Justification of MLE for regression problem

For additive noise models $\mathbf{y} = \mathbf{H}\mathbf{x} + \boldsymbol{\eta}$ where \mathbf{H} is known precisely (all uncertainty is in $\boldsymbol{\eta}$), we have

- ▶ $\operatorname{argmin}_{\mathbf{x}} \|\mathbf{H}\mathbf{x} - \mathbf{y}\|_2^2$ is the MLE for Gaussian noise (ordinary least squares)

Justification of MLE for regression problem

For additive noise models $\mathbf{y} = \mathbf{H}\mathbf{x} + \boldsymbol{\eta}$ where \mathbf{H} is known precisely (all uncertainty is in $\boldsymbol{\eta}$), we have

- ▶ $\operatorname{argmin}_{\mathbf{x}} \|\mathbf{H}\mathbf{x} - \mathbf{y}\|_2^2$ is the MLE for Gaussian noise (ordinary least squares)
- ▶ $\operatorname{argmin}_{\mathbf{x}} \|\mathbf{H}\mathbf{x} - \mathbf{y}\|_1$ is the MLE for double exponential noise

Justification of MLE for regression problem

For additive noise models $\mathbf{y} = \mathbf{H}\mathbf{x} + \boldsymbol{\eta}$ where \mathbf{H} is known precisely (all uncertainty is in $\boldsymbol{\eta}$), we have

- ▶ $\operatorname{argmin}_{\mathbf{x}} \|\mathbf{H}\mathbf{x} - \mathbf{y}\|_2^2$ is the MLE for Gaussian noise (ordinary least squares)
- ▶ $\operatorname{argmin}_{\mathbf{x}} \|\mathbf{H}\mathbf{x} - \mathbf{y}\|_1$ is the MLE for double exponential noise
- ▶ $\operatorname{argmin}_{\mathbf{x}} \|\mathbf{H}\mathbf{x} - \mathbf{y}\|_\infty$ is the MLE for uniform noise.

Justification of MLE for regression problem

For additive noise models $\mathbf{y} = \mathbf{H}\mathbf{x} + \boldsymbol{\eta}$ where \mathbf{H} is known precisely (all uncertainty is in $\boldsymbol{\eta}$), we have

- ▶ $\operatorname{argmin}_{\mathbf{x}} \|\mathbf{H}\mathbf{x} - \mathbf{y}\|_2^2$ is the MLE for Gaussian noise (ordinary least squares)
- ▶ $\operatorname{argmin}_{\mathbf{x}} \|\mathbf{H}\mathbf{x} - \mathbf{y}\|_1$ is the MLE for double exponential noise
- ▶ $\operatorname{argmin}_{\mathbf{x}} \|\mathbf{H}\mathbf{x} - \mathbf{y}\|_\infty$ is the MLE for uniform noise.
- ▶ For uncertainty in operator, total least squares, i.e.,

$$\begin{aligned} & \min_{\mathbf{x}, \mathbf{U}, \boldsymbol{\eta}} \quad \|\mathbf{U}, \boldsymbol{\eta}\|_F \\ & \text{Subject to} \quad (\mathbf{H} + \mathbf{U})\mathbf{x} = \mathbf{y} + \boldsymbol{\eta} \end{aligned}$$

is the MLE for i.i.d. Gaussian in \mathbf{H} and $\boldsymbol{\eta}$

Justification of MLE for regression problem

For additive noise models $\mathbf{y} = \mathbf{H}\mathbf{x} + \boldsymbol{\eta}$ where \mathbf{H} is known precisely (all uncertainty is in $\boldsymbol{\eta}$), we have

- ▶ $\operatorname{argmin}_{\mathbf{x}} \|\mathbf{H}\mathbf{x} - \mathbf{y}\|_2^2$ is the MLE for Gaussian noise (ordinary least squares)
- ▶ $\operatorname{argmin}_{\mathbf{x}} \|\mathbf{H}\mathbf{x} - \mathbf{y}\|_1$ is the MLE for double exponential noise
- ▶ $\operatorname{argmin}_{\mathbf{x}} \|\mathbf{H}\mathbf{x} - \mathbf{y}\|_\infty$ is the MLE for uniform noise.
- ▶ For uncertainty in operator, total least squares, i.e.,

$$\begin{aligned} & \min_{\mathbf{x}, \mathbf{U}, \boldsymbol{\eta}} \quad \|\mathbf{U}, \boldsymbol{\eta}\|_F \\ & \text{Subject to} \quad (\mathbf{H} + \mathbf{U})\mathbf{x} = \mathbf{y} + \boldsymbol{\eta} \end{aligned}$$

is the MLE for i.i.d. Gaussian in \mathbf{H} and $\boldsymbol{\eta}$

MLE for uncertainty in design matrix

- ▶ We seek a PDF for $\mathbf{y} = \mathbf{G}\mathbf{x} + \boldsymbol{\eta}$ to form a likelihood function

MLE for uncertainty in design matrix

- ▶ We seek a PDF for $\mathbf{y} = \mathbf{G}\mathbf{x} + \boldsymbol{\eta}$ to form a likelihood function
- ▶ Each component can be rewritten as a sum of RVs, i.e.,
$$y_i = \sum_{j=1}^n G_{ij}x_j + \eta_i$$

MLE for uncertainty in design matrix

- ▶ We seek a PDF for $\mathbf{y} = \mathbf{G}\mathbf{x} + \boldsymbol{\eta}$ to form a likelihood function
- ▶ Each component can be rewritten as a sum of RVs, i.e.,
$$y_i = \sum_{j=1}^n G_{ij}x_j + \eta_i$$
- ▶ PDFs for sums of RVs are challenging to derive since they generally require convolutions

MLE for uncertainty in design matrix

- ▶ We seek a PDF for $\mathbf{y} = \mathbf{G}\mathbf{x} + \boldsymbol{\eta}$ to form a likelihood function
- ▶ Each component can be rewritten as a sum of RVs, i.e.,
$$y_i = \sum_{j=1}^n G_{ij}x_j + \eta_i$$
- ▶ PDFs for sums of RVs are challenging to derive since they generally require convolutions

Example

Let $Z \sim \mathcal{N}(0, 1)$ and $U \sim \text{Uniform}(0, 1)$. The PDF for $U + Z$ is

$$f_{U+Z}(t) = \frac{1}{\sqrt{2\pi}} \int_0^1 e^{-(t-s)^2/2} ds.$$

No analytic form!

Table of Contents

Problem Formulation

Maximum Likelihood Estimation and Its Limitations

Moment Generating Functions and Saddle Point Approximation

Approximate Likelihood Function and Optimization Problem

Algorithm and Numerical Experiments

Moment generating functions to the rescue

- ▶ The moment generating function (MGF) is a bilateral Laplace transform of PDF given by

$$M_Y(t) = \mathbb{E}(e^{tY})$$

Moment generating functions to the rescue

- ▶ The moment generating function (MGF) is a bilateral Laplace transform of PDF given by

$$M_Y(t) = \mathbb{E}(e^{tY})$$

- ▶ Useful properties, i.e., for indep. RVs U , Z and $a \in \mathbb{R}$

$$M_{aU+Z}(t) = M_U(at)M_Z(t)$$

Moment generating functions to the rescue

- ▶ The moment generating function (MGF) is a bilateral Laplace transform of PDF given by

$$M_Y(t) = \mathbb{E}(e^{tY})$$

- ▶ Useful properties, i.e., for indep. RVs U , Z and $a \in \mathbb{R}$

$$M_{aU+Z}(t) = M_U(at)M_Z(t)$$

- ▶ When MGF exists, it uniquely characterizes the distribution

Moment generating functions to the rescue

- ▶ The moment generating function (MGF) is a bilateral Laplace transform of PDF given by

$$M_Y(t) = \mathbb{E}(e^{tY})$$

- ▶ Useful properties, i.e., for indep. RVs U , Z and $a \in \mathbb{R}$

$$M_{aU+Z}(t) = M_U(at)M_Z(t)$$

- ▶ When MGF exists, it uniquely characterizes the distribution
- ▶ With MGFs, convolution \rightarrow multiplication

Moment generating functions to the rescue

- ▶ The moment generating function (MGF) is a bilateral Laplace transform of PDF given by

$$M_Y(t) = \mathbb{E}(e^{tY})$$

- ▶ Useful properties, i.e., for indep. RVs U , Z and $a \in \mathbb{R}$

$$M_{aU+Z}(t) = M_U(at)M_Z(t)$$

- ▶ When MGF exists, it uniquely characterizes the distribution
- ▶ With MGFs, convolution \rightarrow multiplication

Example

Let $Z \sim \mathcal{N}(0, 1)$ and $U \sim \text{Uniform}(0, 1)$. The MGF for $U + Z$ is

$$M_{U+Z}(t) = \frac{(e^t - 1)e^{-t^2/2}}{t}.$$

Analytic form, no need for quadrature

Moment and generating functions

- ▶ Using properties of MGFs, we have

$$M_{Y_i}(t) = M_{\mathbf{g}_i^T \mathbf{x} + \eta_i}(t) = M_{\eta_i}(t) \prod_{j=1}^n M_{G_{ij}}(tx_j)$$

Moment and generating functions

- ▶ Using properties of MGFs, we have

$$M_{Y_i}(t) = M_{\mathbf{g}_i^T \mathbf{x} + \eta_i}(t) = M_{\eta_i}(t) \prod_{j=1}^n M_{G_{ij}}(tx_j)$$

- ▶ Importantly, we have expressed a complicated MGF for Y_i as the product of simple univariate MGFs for G_{ij} and η_i .

Moment and generating functions

- ▶ Using properties of MGFs, we have

$$M_{Y_i}(t) = M_{\mathbf{g}_i^T \mathbf{x} + \eta_i}(t) = M_{\eta_i}(t) \prod_{j=1}^n M_{G_{ij}}(tx_j)$$

- ▶ Importantly, we have expressed a complicated MGF for Y_i as the product of simple univariate MGFs for G_{ij} and η_i .
- ▶ By inverting transform for M_{Y_i} , we can recover density, but difficult in practice

Moment and generating functions

- ▶ Using properties of MGFs, we have

$$M_{Y_i}(t) = M_{\mathbf{g}_i^T \mathbf{x} + \eta_i}(t) = M_{\eta_i}(t) \prod_{j=1}^n M_{G_{ij}}(tx_j)$$

- ▶ Importantly, we have expressed a complicated MGF for Y_i as the product of simple univariate MGFs for G_{ij} and η_i .
- ▶ By inverting transform for M_{Y_i} , we can recover density, but difficult in practice
- ▶ Use approximation method instead!

Density approximation

Using MGFs to approximate PDFs allows for construction of a likelihood. Some options are

- ▶ Edgeworth series [1]: poor tail behavior (polynomial series)
- ▶ Kernel density estimation [2, 3]: data intensive
- ▶ Saddle point approximation [4, 5, 6]: uses exponential tilting and works well in practice

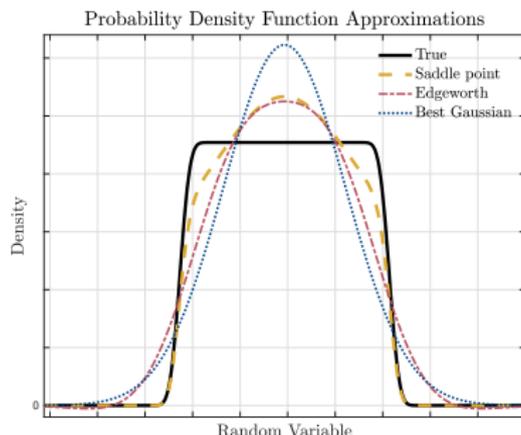


Figure: True density and several approximations when $y = \mathbf{g}^T \mathbf{x} + \eta$

Saddle point approximation

Saddle point approximation for PDF of RV Y is

$$f_Y(y) \approx \sqrt{\frac{1}{2\pi K_Y''(t_0)}} e^{K_Y(t_0) - yt_0}$$

- ▶ $K_Y(t) = \ln M_Y(t)$ is *Cumulant Generating Function* (CGF)
- ▶ t_0 is the solution to $K_Y'(t) - y = 0$ (use Newton's method)

Table of Contents

Problem Formulation

Maximum Likelihood Estimation and Its Limitations

Moment Generating Functions and Saddle Point Approximation

Approximate Likelihood Function and Optimization Problem

Algorithm and Numerical Experiments

Approximate MLE

Using the saddle point approximation and eliminating constants allows us to write the approximate log-likelihood function as

$$\begin{aligned}\ell(\mathbf{x}) &= \sum_{i=1}^m \ln \left\{ \sqrt{\frac{1}{K''_{Y_i}(t_i)}} e^{K_{Y_i}(t_i) - y_i t_i} \right\} \\ &= \sum_{i=1}^m \left\{ K_{Y_i}(t_i) - t_i y_i - \frac{1}{2} \ln (K''_{Y_i}(t_i)) \right\} \\ &= \sum_{i=1}^m \left[K_{\mathbf{g}_i^T \mathbf{x} + \eta_i}(t_i(\mathbf{x})) - \frac{1}{2} \ln (K''_{\mathbf{g}_i^T \mathbf{x} + \eta_i}(t_i(\mathbf{x}))) - t_i(\mathbf{x}) y_i \right].\end{aligned}$$

where \mathbf{g}_i^T is i^{th} row of \mathbf{G} and $t_i(\mathbf{x})$ is solution to $K'_{\mathbf{g}_i^T \mathbf{x} + \eta_i}(t) = y_i$.

Optimization problem

The approximate MLE can be cast generically in vector form

$$\begin{aligned} \operatorname{argmax}_{\mathbf{x}, \mathbf{t}} \quad & \mathbf{1}^T \left(K_{\mathbf{G}\mathbf{x}+\boldsymbol{\eta}}(\mathbf{t}) - \frac{1}{2} \ln (K''_{\mathbf{G}\mathbf{x}+\boldsymbol{\eta}}(\mathbf{t})) \right) - \mathbf{t}^T \mathbf{y} \\ \text{Subject to} \quad & K'_{\mathbf{G}\mathbf{x}+\boldsymbol{\eta}}(\mathbf{t}) = \mathbf{y} \end{aligned}$$

Optimization problem

The approximate MLE can be cast generically in vector form

$$\begin{aligned} \operatorname{argmax}_{\mathbf{x}, \mathbf{t}} \quad & \mathbf{1}^T \left(K_{\mathbf{G}\mathbf{x}+\boldsymbol{\eta}}(\mathbf{t}) - \frac{1}{2} \ln (K''_{\mathbf{G}\mathbf{x}+\boldsymbol{\eta}}(\mathbf{t})) \right) - \mathbf{t}^T \mathbf{y} \\ \text{Subject to} \quad & K'_{\mathbf{G}\mathbf{x}+\boldsymbol{\eta}}(\mathbf{t}) = \mathbf{y} \end{aligned}$$

Example

When $\mathbf{G} \sim \text{Uniform}(\mathbf{H} - \delta \mathbf{1} \mathbf{1}^T, \mathbf{H} + \delta \mathbf{1} \mathbf{1}^T)$ and $\boldsymbol{\eta} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$,

$$\begin{aligned} \operatorname{argmax}_{\mathbf{x}, \mathbf{t}} \quad & \mathbf{t}^T \left(\frac{\sigma^2}{2} \mathbf{t} + \mathbf{H} \mathbf{x} - \mathbf{y} \right) + \mathbf{1}^T \ln [\sinh (\delta \mathbf{t} \mathbf{x}^T) \circ (\delta \mathbf{t} \mathbf{x}^T)] \mathbf{1} \\ & - \frac{1}{2} \mathbf{1}^T \ln [\sigma^2 \mathbf{1} - \delta^2 \operatorname{csch}^2 (\delta \mathbf{t} \mathbf{x}^T) \mathbf{x}^2] \end{aligned}$$

$$\text{Subject to} \quad \sigma^2 \mathbf{t} + \mathbf{H} \mathbf{x} + \delta \operatorname{coth} (\delta \mathbf{t} \mathbf{x}^T) \mathbf{x} - n(\mathbf{1} \circ \mathbf{t}) = \mathbf{y}.$$

Gradients

- ▶ Despite constraint requiring a numeric solve, have gradients

Gradients

- ▶ Despite constraint requiring a numeric solve, have gradients
- ▶ Letting $\mathbf{q}(\mathbf{x}, \mathbf{t}) = K'_{\mathbf{G}\mathbf{x}+\eta}(\mathbf{t}) - \mathbf{y}$ be our constraint, then using adjoint state method [7], we have

$$\nabla_{\mathbf{x}}\ell = \frac{\partial\ell}{\partial\mathbf{x}} - \left(\frac{\partial\ell}{\partial\mathbf{t}}\right) \left(\frac{\partial\mathbf{q}}{\partial\mathbf{t}}\right)^{-1} \left(\frac{\partial\mathbf{q}}{\partial\mathbf{x}}\right). \quad (1)$$

Gradients

- ▶ Despite constraint requiring a numeric solve, have gradients
- ▶ Letting $\mathbf{q}(\mathbf{x}, \mathbf{t}) = K'_{\mathbf{G}_{\mathbf{x}+\eta}}(\mathbf{t}) - \mathbf{y}$ be our constraint, then using adjoint state method [7], we have

$$\nabla_{\mathbf{x}} \ell = \frac{\partial \ell}{\partial \mathbf{x}} - \left(\frac{\partial \ell}{\partial \mathbf{t}} \right) \left(\frac{\partial \mathbf{q}}{\partial \mathbf{t}} \right)^{-1} \left(\frac{\partial \mathbf{q}}{\partial \mathbf{x}} \right). \quad (1)$$

Unimportant, but for completeness, each factor is given by:

$$\frac{\partial \ell}{\partial \mathbf{x}} = \mathbf{1}^T \left(\frac{\partial}{\partial \mathbf{x}} K_{\mathbf{G}_{\mathbf{x}+\eta}}(\mathbf{t}) - \frac{1}{2} \{ \text{diag} (K''_{\mathbf{G}_{\mathbf{x}+\eta}}(\mathbf{t})) \}^{-1} \frac{\partial}{\partial \mathbf{x}} K''_{\mathbf{G}_{\mathbf{x}+\eta}}(\mathbf{t}) \right)$$

$$\frac{\partial \ell}{\partial \mathbf{t}} = \left(K'_{\mathbf{G}_{\mathbf{x}+\eta}}(\mathbf{t}) - \frac{1}{2} (K'''_{\mathbf{G}_{\mathbf{x}+\eta}}(\mathbf{t}) \oslash K''_{\mathbf{G}_{\mathbf{x}+\eta}}(\mathbf{t})) - \mathbf{y} \right)^T$$

$$\frac{\partial \mathbf{q}}{\partial \mathbf{t}} = \text{diag} (K''_{\mathbf{G}_{\mathbf{x}+\eta}}(\mathbf{t})) ,$$

$$\frac{\partial \mathbf{q}}{\partial \mathbf{x}} = \frac{\partial}{\partial \mathbf{x}} K'_{\mathbf{G}_{\mathbf{x}+\eta}}(\mathbf{t}).$$

Table of Contents

Problem Formulation

Maximum Likelihood Estimation and Its Limitations

Moment Generating Functions and Saddle Point Approximation

Approximate Likelihood Function and Optimization Problem

Algorithm and Numerical Experiments

- ▶ With gradient information, we can employ first order methods

Algorithms and Numerical Experiments

- ▶ With gradient information, we can employ first order methods
- ▶ For experiments, we opted for L-BFGS [8], a quasi-Newton method, which solved problem rapidly

Algorithms and Numerical Experiments

- ▶ With gradient information, we can employ first order methods
- ▶ For experiments, we opted for L-BFGS [8], a quasi-Newton method, which solved problem rapidly
- ▶ Although possible to calculate derivatives analytically in many cases, automatic differentiation can save time and trouble [9]

Numerical Experiments

- ▶ To validate our method, we solved 10,000 problems. Using the generative model $\mathbf{y} = \mathbf{G}\mathbf{x}_{\text{TRU}} + \boldsymbol{\eta}$ with components drawn as follows:

Numerical Experiments

- ▶ To validate our method, we solved 10,000 problems. Using the generative model $\mathbf{y} = \mathbf{G}\mathbf{x}_{\text{TRU}} + \boldsymbol{\eta}$ with components drawn as follows:
 - ▶ $\mathbf{G} \in \mathbb{R}^{m \times n}$, from continuous uniform matrices on $[0, 10]$

Numerical Experiments

- ▶ To validate our method, we solved 10,000 problems. Using the generative model $\mathbf{y} = \mathbf{G}\mathbf{x}_{\text{TRU}} + \boldsymbol{\eta}$ with components drawn as follows:
 - ▶ $\mathbf{G} \in \mathbb{R}^{m \times n}$, from continuous uniform matrices on $[0, 10]$
 - ▶ $\mathbf{x}_{\text{TRU}} \in \mathbb{R}^n$, taken from the heavy-tailed Cauchy distribution to make the use of prior information on the solution difficult

Numerical Experiments

- ▶ To validate our method, we solved 10,000 problems. Using the generative model $\mathbf{y} = \mathbf{G}\mathbf{x}_{\text{TRU}} + \boldsymbol{\eta}$ with components drawn as follows:
 - ▶ $\mathbf{G} \in \mathbb{R}^{m \times n}$, from continuous uniform matrices on $[0, 10]$
 - ▶ $\mathbf{x}_{\text{TRU}} \in \mathbb{R}^n$, taken from the heavy-tailed Cauchy distribution to make the use of prior information on the solution difficult
 - ▶ $\boldsymbol{\eta} \in \mathbb{R}^m$ has i.i.d. components drawn from a normal distribution, i.e., $\eta_i \sim \mathcal{N}(0, \sigma^2)$

Numerical Experiments

- ▶ To validate our method, we solved 10,000 problems. Using the generative model $\mathbf{y} = \mathbf{G}\mathbf{x}_{\text{TRU}} + \boldsymbol{\eta}$ with components drawn as follows:
 - ▶ $\mathbf{G} \in \mathbb{R}^{m \times n}$, from continuous uniform matrices on $[0, 10]$
 - ▶ $\mathbf{x}_{\text{TRU}} \in \mathbb{R}^n$, taken from the heavy-tailed Cauchy distribution to make the use of prior information on the solution difficult
 - ▶ $\boldsymbol{\eta} \in \mathbb{R}^m$ has i.i.d. components drawn from a normal distribution, i.e., $\eta_i \sim \mathcal{N}(0, \sigma^2)$
- ▶ Assuming knowledge of σ^2 , and inferred value of δ by observing $\mathbf{H} = \text{round}(\mathbf{G})$ and \mathbf{y}

Numerical Experiments

- ▶ To validate our method, we solved 10,000 problems. Using the generative model $\mathbf{y} = \mathbf{G}\mathbf{x}_{\text{TRU}} + \boldsymbol{\eta}$ with components drawn as follows:
 - ▶ $\mathbf{G} \in \mathbb{R}^{m \times n}$, from continuous uniform matrices on $[0, 10]$
 - ▶ $\mathbf{x}_{\text{TRU}} \in \mathbb{R}^n$, taken from the heavy-tailed Cauchy distribution to make the use of prior information on the solution difficult
 - ▶ $\boldsymbol{\eta} \in \mathbb{R}^m$ has i.i.d. components drawn from a normal distribution, i.e., $\eta_i \sim \mathcal{N}(0, \sigma^2)$
- ▶ Assuming knowledge of σ^2 , and inferred value of δ by observing $\mathbf{H} = \text{round}(\mathbf{G})$ and \mathbf{y}
- ▶ Used the proposed approximate MLE to estimate \mathbf{x}

Numerical Experiments

- ▶ To validate our method, we solved 10,000 problems. Using the generative model $\mathbf{y} = \mathbf{G}\mathbf{x}_{\text{TRU}} + \boldsymbol{\eta}$ with components drawn as follows:
 - ▶ $\mathbf{G} \in \mathbb{R}^{m \times n}$, from continuous uniform matrices on $[0, 10]$
 - ▶ $\mathbf{x}_{\text{TRU}} \in \mathbb{R}^n$, taken from the heavy-tailed Cauchy distribution to make the use of prior information on the solution difficult
 - ▶ $\boldsymbol{\eta} \in \mathbb{R}^m$ has i.i.d. components drawn from a normal distribution, i.e., $\eta_i \sim \mathcal{N}(0, \sigma^2)$
- ▶ Assuming knowledge of σ^2 , and inferred value of δ by observing $\mathbf{H} = \text{round}(\mathbf{G})$ and \mathbf{y}
- ▶ Used the proposed approximate MLE to estimate \mathbf{x}
- ▶ Compared to ordinary least squares and total least squares

Numerical Experiments

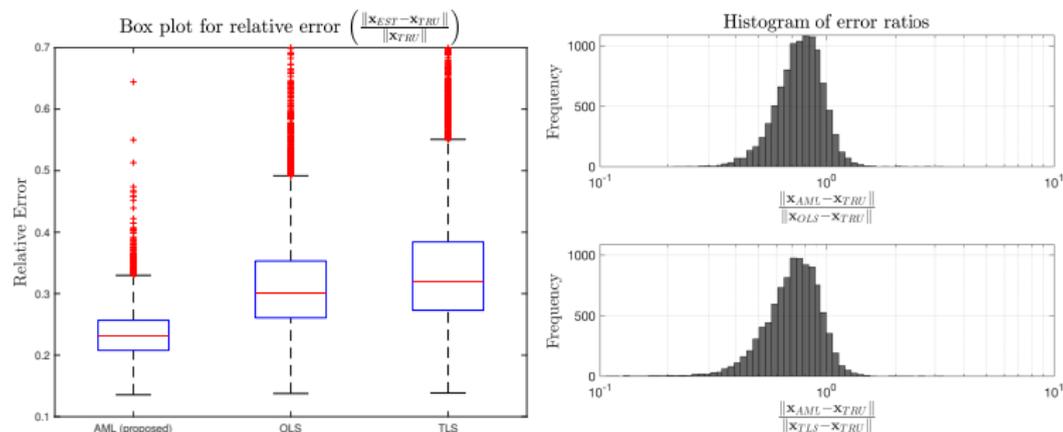


Figure: Error metrics for simulations $\mathbf{G} \in \mathbb{R}^{110 \times 100}$ and $\sigma = 1$ over 10,000 simulations. Design matrix rounded to ones spot. Left: box-plot of relative error for different methods. Right: histogram of error ratio $\frac{\|\mathbf{x}_{AML} - \mathbf{x}_{TRU}\|}{\|\mathbf{x}_{OLS} - \mathbf{x}_{TRU}\|}$. Values less than 1 indicate AML outperformed competing method for identical data.

Summary

- ▶ Maximum likelihood estimation is useful, but forming likelihood functions for general noise models is difficult

Summary

- ▶ Maximum likelihood estimation is useful, but forming likelihood functions for general noise models is difficult
- ▶ Presented a method to construct an approximate likelihood function based on MGFs and the saddle point approximation to avoid difficulties

Summary

- ▶ Maximum likelihood estimation is useful, but forming likelihood functions for general noise models is difficult
- ▶ Presented a method to construct an approximate likelihood function based on MGFs and the saddle point approximation to avoid difficulties
- ▶ Found gradient of approximate likelihood using the adjoint state method allowing use of off-the-shelf algorithms

Summary

- ▶ Maximum likelihood estimation is useful, but forming likelihood functions for general noise models is difficult
- ▶ Presented a method to construct an approximate likelihood function based on MGFs and the saddle point approximation to avoid difficulties
- ▶ Found gradient of approximate likelihood using the adjoint state method allowing use of off-the-shelf algorithms
- ▶ Showed results of numerical experiments illustrating its effectiveness

References

- [1] P. Hall, *The bootstrap and Edgeworth expansion*. Springer Science & Business Media, 2013.
- [2] E. Parzen, "On estimation of a probability density function and mode," *The annals of mathematical statistics*, vol. 33, no. 3, pp. 1065–1076, 1962.
- [3] M. Rosenblatt, "Remarks on some nonparametric estimates of a density function," *The annals of mathematical statistics*, vol. 27, no. 3, pp. 832–837, 1956.
- [4] H. E. Daniels, "Saddlepoint approximations in statistics," *The annals of mathematical statistics*, pp. 631–650, 1954.
- [5] R. Lugannani and S. Rice, "Saddle point approximation for the distribution of the sum of independent random variables," *Advances in applied probability*, vol. 12, no. 2, pp. 475–490, 1980.
- [6] O. Barndorff-Nielsen and D. R. Cox, "Edgeworth and saddle-point approximations with statistical applications," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 41, no. 3, pp. 279–299, 1979.
- [7] G. Strang, *Computational methods for inverse problems*, vol. 23. Wellesley-Cambridge Press, 2007.
- [8] J. Nocedal, "Updating quasi-Newton matrices with limited storage," *Math. Comp.*, vol. 25, no. 151, pp. 773–782, 1980.
- [9] M. J. Weinstein and A. V. Rao, "Algorithm 984: Adigator, a toolbox for the algorithmic differentiation of mathematical functions in matlab using source transformation via operator overloading," *ACM Transactions on Mathematical Software (TOMS)*, vol. 44, no. 2, p. 21, 2017.

Thank you for your time!

Questions?